

# Approximate Dynamic Programming via a Smoothed Linear Program

Vijay V. Desai

Industrial Engineering and Operations Research

Columbia University

email: [vvd2101@columbia.edu](mailto:vvd2101@columbia.edu)

Vivek F. Farias

Sloan School of Management

Massachusetts Institute of Technology

email: [vivekf@mit.edu](mailto:vivekf@mit.edu)

Ciamac C. Moallemi

Graduate School of Business

Columbia University

email: [ciamac@gsb.columbia.edu](mailto:ciamac@gsb.columbia.edu)

Initial Version: August 4, 2009

Current Revision: September 25, 2009

## Abstract

We present a novel linear program for the approximation of the dynamic programming cost-to-go function in high-dimensional stochastic control problems. LP approaches to approximate DP have typically relied on a natural ‘projection’ of a well studied linear program for exact dynamic programming. Such programs restrict attention to approximations that are lower bounds to the optimal cost-to-go function. Our program—the ‘smoothed approximate linear program’—is distinct from such approaches and relaxes the restriction to lower bounding approximations in an appropriate fashion while remaining computationally tractable. Doing so appears to have several advantages: First, we demonstrate substantially superior bounds on the quality of approximation to the optimal cost-to-go function afforded by our approach. Second, experiments with our approach on a challenging problem (the game of Tetris) show that the approach outperforms the existing LP approach (which has previously been shown to be competitive with several ADP algorithms) by an order of magnitude.

## 1. Introduction

Many dynamic optimization problems can be cast as Markov decision problems (MDPs) and solved, in principle, via dynamic programming. Unfortunately, this approach is frequently untenable due to the ‘curse of dimensionality’. Approximate dynamic programming (ADP) is an approach which attempts to address this difficulty. ADP algorithms seek to compute good approximations to the dynamic programming optimal cost-to-go function within the span of some pre-specified set of basis functions.

ADP algorithms are typically motivated by exact algorithms for dynamic programming. The approximate linear programming (ALP) method is one such approach, motivated by the LP used

for the computation of the optimal cost-to-go function. Introduced by Schweitzer and Seidmann (1985) and analyzed and further developed by de Farias and Van Roy (2003, 2004), this approach is attractive for a number of reasons. First, the availability of efficient solvers for linear programming makes the LP approach easy to implement. Second, the approach offers attractive theoretical guarantees. In particular, the quality of the approximation to the cost-to-go function produced by the LP approach can be shown to compete, in an appropriate sense, with the quality of the best possible approximation afforded by the set of basis functions used. A testament to the success of the LP approach is the number of applications it has seen in recent years in large scale dynamic optimization problems. These applications range from the control of queueing networks to revenue management to the solution of large scale stochastic games.

The optimization program employed in the ALP approach is in some sense the most natural linear programming formulation for ADP. In particular, the ALP is identical to the linear program used for exact computation of the optimal cost-to-go function, with further constraints limiting solutions to the low-dimensional subspace spanned by the basis functions used. The resulting LP implicitly restricts attention to approximations that are lower bounds to the optimal cost-to-go function. The structure of this program appears crucial in establishing guarantees on the quality of approximations produced by the approach; these approximation guarantees were remarkable and a first for any ADP method. That said, the restriction to lower bounds naturally leads one to ask whether the program employed by the ALP approach is the ‘right’ math programming formulation for ADP. In particular, it may be advantageous to relax the lower bound requirement so as to allow for a better approximation, and, ultimately, better policy performance. Is there an alternative formulation that permits better approximations to the cost-to-go function while remaining computationally tractable? Motivated by this question, the present paper introduces a new linear program for ADP we call the ‘smoothed’ approximate linear program (or SALP). We believe that the SALP provides a preferable math programming formulation for ADP. In particular, we make the following contributions:

1. We are able to establish strong approximation and performance guarantees for approximations to the cost-to-go function produced by the SALP; these guarantees are *substantially* stronger than the corresponding guarantees for the ALP.
2. The number of constraints and variables in the SALP scale with the size of the MDP state space. We nonetheless establish sample complexity bounds that demonstrate that an appropriate ‘sampled’ SALP provides a good approximation to the SALP solution with a tractable number of sampled MDP states. Moreover, we identify structural properties for the sampled SALP that can be exploited for fast optimization. Our sample complexity results and these structural observations allow us to conclude that the SALP is essentially no harder to solve than existing LP formulations for ADP.
3. We present a computational study demonstrating the efficacy of our approach on the game

of Tetris. Tetris is a notoriously difficult, ‘unstructured’ dynamic optimization problem and has been used as a convenient testbed problem for numerous ADP approaches. The ALP has been demonstrated to be competitive with other ADP approaches for Tetris, such as temporal difference learning or policy gradient methods (see Farias and Van Roy, 2006). In detailed comparisons with the ALP, we show that the SALP provides an *order of magnitude* improvement over controllers designed via that approach for the game of Tetris.

The literature on ADP algorithms is vast and we make no attempt to survey it here. Van Roy (2002) or Bertsekas (2007, Chap. 6) provide good, brief overviews, while Bertsekas and Tsitsiklis (1996) and Powell (2007) are encyclopedic references on the topic. The exact LP for the solution of dynamic programs is attributed to Manne (1960). The ALP approach to ADP was introduced by Schweitzer and Seidmann (1985) and de Farias and Van Roy (2003, 2004). de Farias and Van Roy (2003) establish strong approximation guarantees for ALP based approximations assuming knowledge of a ‘Lyapunov’-like function which must be included in the basis. The approach we present may be viewed as optimizing over all possible Lyapunov functions. de Farias and Van Roy (2006) introduce a program for average cost approximate dynamic programming that resembles the SALP; a critical difference is that their program requires the relative violation allowed across ALP constraints be specified as input. Applications of the LP approach to ADP range from scheduling in queueing networks (Morrison and Kumar., 1999; Veatch, 2005; Moallemi et al., 2008), revenue management (Adelman, 2007; Farias and Van Roy, 2007; Zhang and Adelman, 2008), portfolio management (Han, 2005), inventory problems (Adelman, 2004; Adelman and Klabjan, 2009), and algorithms for solving stochastic games (Farias et al., 2008) among others. Remarkably, in applications such as network revenue management, control policies produced via the LP approach (namely, Adelman, 2007; Farias and Van Roy, 2007) are competitive with ADP approaches that carefully exploit problem structure, such as for instance that of Topaloglu (2009).

The remainder of this paper is organized as follows: In Section 2, we formulate the approximate dynamic programming setting and describe the ALP approach. The smoothed ALP is developed as a relaxation of the ALP in Section 3. Section 4 provides a theoretical analysis of the SALP, in terms of approximation and performance guarantees, as well as a sample complexity bound. In Section 5, we describe the practical implementation of the SALP method, illustrating how parameter choices can be made as well as how to efficiently solve the resulting optimization program. Section 6 contains the computational study of the game Tetris. Finally, in Section 7, we conclude.

## 2. Problem Formulation

Our setting is that of a discrete-time, discounted infinite-horizon, cost-minimizing MDP with a finite state space  $\mathcal{X}$  and finite action space  $\mathcal{A}$ . At time  $t$ , given the current state  $x_t$  and a choice of action  $a_t$ , a per-stage cost  $g(x_t, a_t)$  is incurred. The subsequent state  $x_{t+1}$  is determined according to the transition probability kernel  $P_{a_t}(x_t, \cdot)$ .

A stationary policy  $\mu: \mathcal{X} \rightarrow \mathcal{A}$  is a mapping that determines the choice of action at each time as a function of the state. Given each initial state  $x_0 = x$ , the expected discounted cost (cost-to-go function) of the policy  $\mu$  is given by

$$J_\mu(x) \triangleq \mathbb{E}_\mu \left[ \sum_{t=0}^{\infty} \alpha^t g(x_t, \mu(x_t)) \mid x_0 = x \right].$$

Here,  $\alpha \in (0, 1)$  is the discount factor. The expectation is taken under the assumption that actions are selected according to the policy  $\mu$ . In other words, at each time  $t$ ,  $a_t \triangleq \mu(x_t)$ .

Denote by  $P_\mu \in \mathbb{R}^{\mathcal{X} \times \mathcal{X}}$  the transition probability matrix for the policy  $\mu$ , whose  $(x, x')$ th entry is  $P_{\mu(x)}(x, x')$ . Denote by  $g_\mu \in \mathbb{R}^{\mathcal{X}}$  the vector whose  $x$ th entry is  $g(x, \mu(x))$ . Then, the cost-to-go function  $J_\mu$  can be written in vector form as

$$J_\mu = \sum_{t=0}^{\infty} \alpha^t P_\mu^t g_\mu.$$

Further, the cost-to-go function  $J_\mu$  is the unique solution to the equation  $T_\mu J = J$ , where the operator  $T_\mu$  is defined by  $T_\mu J = g_\mu + \alpha P_\mu J$ .

Our goal is to find an optimal stationary policy  $\mu^*$ , that is, a policy that minimizes the expected discounted cost from every state  $x$ . In particular,

$$\mu^* \in \underset{\mu}{\operatorname{argmin}} J_\mu(x).$$

The Bellman operator  $T$  is defined component-wise according to

$$(TJ)(x) \triangleq \min_{a \in \mathcal{A}} g(x, a) + \alpha \sum_{x' \in \mathcal{X}} P_a(x, x') J(x'), \quad \forall x \in \mathcal{X}.$$

Bellman's equation is then the fixed point equation

$$(1) \quad TJ = J.$$

Standard results in dynamic programming establish that the optimal cost-to-go function  $J^*$  is the unique solution to Bellman's equation (see, for example, Bertsekas, 2007, Chap. 1). Further, if  $\mu^*$  is a policy that is greedy with respect to  $J^*$  (i.e.,  $\mu^*$  satisfies  $TJ^* = T_{\mu^*}J^*$ ), then  $\mu^*$  is an optimal policy.

## 2.1. The Linear Programming Approach

A number of computational approaches are available for the solution of the Bellman equation. One approach involves solving the optimization program:

$$(2) \quad \begin{aligned} & \underset{J}{\text{maximize}} && \nu^\top J \\ & \text{subject to} && J \leq TJ. \end{aligned}$$

Here,  $\nu \in \mathbb{R}^{\mathcal{X}}$  is a vector with positive components that are known as the *state-relevance weights*. The above program is indeed an LP since for each state  $x$ , the constraint  $J(x) \leq (TJ)(x)$  is equivalent to the set of  $|\mathcal{A}|$  linear constraints

$$J(x) \leq g(x, a) + \alpha \sum_{x' \in \mathcal{X}} P_a(x, x') J(x'), \quad \forall a \in \mathcal{A}.$$

We refer to (2) as the *exact LP*.

Suppose that a vector  $J$  is feasible for the exact LP (2). Since  $J \leq TJ$ , monotonicity of the Bellman operator implies that  $J \leq T^k J$ , for any integer  $k \geq 1$ . Since the Bellman operator  $T$  is a contraction,  $T^k J$  must converge to the unique fixed point  $J^*$  as  $k \rightarrow \infty$ . Thus, we have that  $J \leq J^*$ . Then, it is clear that every feasible point for (2) is a component-wise lower bound to  $J^*$ . Since  $J^*$  itself is feasible for (2), it must be that  $J^*$  is the unique optimal solution to the exact LP.

## 2.2. The Approximate Linear Program

In many problems, the size of the state space is enormous due to the curse of dimensionality. In such cases, it may be prohibitive to store, much less compute, the optimal cost-to-go function  $J^*$ . In approximate dynamic programming (ADP), the goal is to find tractable approximations to the optimal cost-to-go function  $J^*$ , with the hope that they will lead to good policies.

Specifically, consider a collection of *basis functions*  $\{\phi_1, \dots, \phi_K\}$  where each  $\phi_i: \mathcal{X} \rightarrow \mathbb{R}$  is a real-valued function on the state space. ADP algorithms seek to find linear combinations of the basis functions that provide good approximations to the optimal cost-to-go function. In particular, we seek a vector of weights  $r \in \mathbb{R}^K$  so that

$$J^*(x) \approx J_r(x) \triangleq \sum_{i=1}^K \phi_i(x) r_i = \Phi r(x).$$

Here, we define  $\Phi \triangleq [\phi_1 \ \phi_2 \ \dots \ \phi_K]$  to be a matrix with columns consisting of the basis functions. Given a vector of weights  $r$  and the corresponding value function approximation  $\Phi r$ , a policy  $\mu_r$  is naturally defined as the ‘greedy’ policy with respect to  $\Phi r$ , i.e. as  $T_{\mu_r} \Phi r = T \Phi r$ .

One way to obtain a set of weights is to solve the exact LP (2), but restricting to the low-dimensional subspace of vectors spanned by the basis functions. This leads to the *approximate*

linear program (ALP), which is defined by

$$(3) \quad \begin{aligned} & \underset{r}{\text{maximize}} && \nu^\top \Phi r \\ & \text{subject to} && \Phi r \leq T\Phi r. \end{aligned}$$

For the balance of the paper, we will make the following assumption:

**Assumption 1.** Assume the  $\nu$  is a probability distribution ( $\nu \geq 0$ ,  $\mathbf{1}^\top \nu = 1$ ), and that the constant function  $\mathbf{1}$  is in the span of the basis functions  $\Phi$ .

The geometric intuition behind the ALP is illustrated in Figure 1(a). Supposed that  $r_{\text{ALP}}$  is a vector that is optimal for the ALP. Then the approximate value function  $\Phi r_{\text{ALP}}$  will lie on the subspace spanned by the columns of  $\Phi$ , as illustrated by the orange line.  $\Phi r_{\text{ALP}}$  will also satisfy the constraints of the exact LP, illustrated by the dark gray region. By the discussion in Section 2.1, this implies that  $\Phi r_{\text{ALP}} \leq J^*$ . In other words, the approximate cost-to-go function is necessarily a point-wise lower bound to the true cost-to-go function in the span of  $\Phi$ .

One can thus interpret the ALP solution  $r_{\text{ALP}}$  equivalently as the optimal solution to the program

$$(4) \quad \begin{aligned} & \underset{r}{\text{minimize}} && \|J^* - \Phi r\|_{1,\nu} \\ & \text{subject to} && \Phi r \leq T\Phi r. \end{aligned}$$

Here, the weighted 1-norm in the objective is defined by

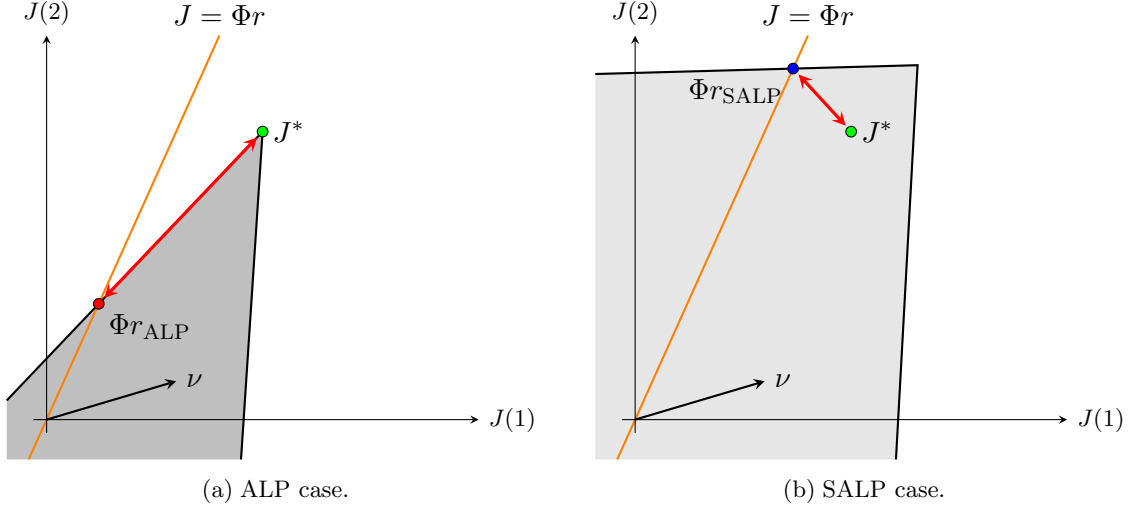
$$\|J^* - \Phi r\|_{1,\nu} \triangleq \sum_{x \in \mathcal{X}} \nu(x) |J^*(x) - \Phi r(x)|.$$

This implies that the approximate LP will find the closest approximation (in the appropriate norm) to the optimal cost-to-go function, out of all approximations satisfying the constraints of the exact LP.

### 3. The Smoothed ALP

The  $J \leq TJ$  constraints in the exact LP, which carry over to the ALP, impose a strong restriction on the cost-to-go function approximation: in particular they restrict us to approximations that are lower bounds to  $J^*$  at *every point in the state space*. In the case where the state space is very large, and the number of basis functions is (relatively) small, it may be the case that constraints arising from rarely visited or pathological states are binding and influence the optimal solution.

In many cases, the ultimate goal is not to find a *lower bound* on the optimal cost-to-go function, but rather to find a *good approximation*. In these instances, it may be that relaxing the constraints in the ALP, so as not to require a uniform lower bound, may allow for better overall approximations to the optimal cost-to-go function. This is also illustrated in Figure 1. Relaxing the feasible region



**Figure 1:** A cartoon illustrating the feasible set and optimal solution for the ALP and SALP, in the case of a two-state MDP. The axes correspond to the components of the value function. A careful relaxation from the feasible set of the ALP to that of the SALP can yield an improved approximation. It is easy to construct a concrete two state example with the above features.

of the ALP in Figure 1(a) to the light gray region in Figure 1(b) would yield the point  $\Phi r_{\text{SALP}}$  as an optimal solution. The relaxation in this case is clearly beneficial; it allows us to compute a better approximation to  $J^*$  than the point  $\Phi r_{\text{SALP}}$ .

Can we construct a fruitful relaxation of this sort in general? The *smoothed approximate linear program* (SALP) is given by:

$$\begin{aligned}
 (5) \quad & \underset{r,s}{\text{maximize}} && \nu^\top \Phi r \\
 & \text{subject to} && \Phi r \leq T\Phi r + s, \\
 & && \pi^\top s \leq \theta, \quad s \geq 0.
 \end{aligned}$$

Here, a vector  $s \in \mathbb{R}^{\mathcal{X}}$  of additional decision variables has been introduced. For each state  $x$ ,  $s(x)$  is a non-negative decision variable (a slack) that allows for violation of the corresponding ALP constraint. The parameter  $\theta \geq 0$  is a non-negative scalar. The parameter  $\pi \in \mathbb{R}^{\mathcal{X}}$  is a probability distribution known as the *constraint violation distribution*. The parameter  $\theta$  is thus a *violation budget*: the expected violation of the  $\Phi r \leq T\Phi r$  constraint, under the distribution  $\pi$ , must be less than  $\theta$ .

The SALP can be alternatively written as

$$\begin{aligned}
 (6) \quad & \underset{r}{\text{maximize}} && \nu^\top \Phi r \\
 & \text{subject to} && \pi^\top (\Phi r - T\Phi r)^+ \leq \theta.
 \end{aligned}$$

Here, given a vector  $J$ ,  $J^+(x) \triangleq \max(J(x), 0)$  is defined to be the component-wise positive part.

Note that, when  $\theta = 0$ , the SALP is equivalent to the ALP. When  $\theta > 0$ , the SALP replaces the ‘hard’ constraints of the ALP with ‘soft’ constraints in the form of a hinge-loss function.

The balance of the paper is concerned with establishing that the SALP forms the basis of a useful approximate dynamic programming algorithm in large scale problems:

- We identify a concrete choice of violation budget  $\theta$  and an idealized constraint violation distribution  $\pi$  for which the SALP provides a useful relaxation in that the optimal solution can be a better approximation to the optimal cost-to-go function. This brings the cartoon improvement in Figure 1 to fruition for general problems.
- We show that the SALP is tractable (i.e., it is well approximated by an appropriate ‘sampled’ version) and present computational experiments for a hard problem (Tetris) illustrating an order of magnitude improvement over the ALP.

## 4. Analysis

This section is dedicated to a theoretical analysis of the SALP. The overarching objective of this analysis is to provide some assurance of the soundness of the proposed approach. In some instances, the bounds we provide will be directly comparable to bounds that have been developed for the ALP method. As such, a relative consideration of the bounds in these two cases can provide a theoretical comparison between the ALP and SALP methods. In addition, our analysis will serve as a crucial guide to practical implementation of the SALP as will be described in Section 5. In particular, the theoretical analysis presented here provides intuition as to how to select parameters such as the state-relevance weights and the constraint violation distribution. We note that all of our bounds are relative to a measure of how well the approximation architecture employed is capable of approximating the optimal cost-to-go function; it is unreasonable to expect non-trivial bounds that are independent of the architecture used.

Our analysis will present three types of results:

- Approximation guarantees (Sections 4.2 and 4.3): We establish bounds on the distance between approximations computed by the SALP and the optimal value function  $J^*$ , relative to the distance between the best possible approximation afforded by the chosen basis functions and  $J^*$ . These guarantees will indicate that the SALP computes approximations that are of comparable quality to the projection<sup>1</sup> of  $J^*$  on to the linear span of  $\Phi$ .
- Performance bounds (Section 4.4): While it is desirable to approximate  $J^*$  as closely as possible, an important concern is the quality of the policies generated by acting greedily according to such approximations, as measured by their performance. We present bounds on the performance loss incurred, relative to the optimal policy, in using an SALP approximation.

---

<sup>1</sup>Note that it is intractable to directly compute the projection since  $J^*$  is unknown.



- Sample complexity results (Section 4.5): The SALP is a linear program with a large number of constraints as well as variables. In practical implementations, one may consider a ‘sampled’ version of this program that has a manageable number of variables and constraints. We present sample complexity guarantees that establish bounds on the number of samples required to produce a good approximation to the solution of the SALP. These bounds scale linearly with the number of basis function  $K$  and are independent of the size of the state space  $\mathcal{X}$ .

#### 4.1. Idealized Assumptions

Given the broad scope of problems addressed by ADP algorithms, analyses of such algorithms typically rely on an ‘idealized’ assumption of some sort. In the case of the ALP, one either assumes the ability to solve a linear program with as many constraints as there are states, or, absent that, knowledge of a certain idealized sampling distribution, so that one can then proceed with solving a ‘sampled’ version of the ALP. Our analysis of the SALP in this section is predicated on the knowledge of this same idealized sampling distribution. In particular, letting  $\mu^*$  be an optimal policy and  $P_{\mu^*}$  the associated transition matrix, we will require access to samples drawn according to the distribution  $\pi_{\mu^*,\nu}$  given by

$$(7) \quad \pi_{\mu^*,\nu}^\top \triangleq (1 - \alpha)\nu^\top(I - \alpha P_{\mu^*})^{-1} = (1 - \alpha) \sum_{t=0}^{\infty} \alpha^t \nu^\top P_{\mu^*}^t.$$

Here  $\nu$  is an arbitrary initial distribution over states. The distribution  $\pi_{\mu^*,\nu}$  may be interpreted as yielding the discounted expected frequency of visits to a given state when the initial state is distributed according to  $\nu$  and the system runs under the policy  $\mu^*$ . We note that the ‘sampled’ ALP introduced by de Farias and Van Roy (2004) requires access to states sampled according to precisely this distribution. Theoretical analyses of other approaches to approximate DP such as approximate value iteration and temporal difference learning similarly rely on the knowledge of specialized sampling distributions that cannot be obtained tractably (see de Farias and Van Roy, 2000).

#### 4.2. A Simple Approximation Guarantee

This section presents a first, simple approximation guarantee for the following specialization of the SALP in (5),

$$(8) \quad \begin{aligned} & \underset{r,s}{\text{maximize}} && \nu^\top \Phi r \\ & \text{subject to} && \Phi r \leq T\Phi r + s, \\ & && \pi_{\mu^*,\nu}^\top s \leq \theta, \quad s \geq 0. \end{aligned}$$

Here, the constraint violation distribution is set to be  $\pi_{\mu^*,\nu}$ .

Before we state our approximation guarantee, consider the following function:

$$(9) \quad \begin{aligned} \ell(r, \theta) &\triangleq \underset{s, \gamma}{\text{minimize}} && \gamma/(1 - \alpha) \\ &\text{subject to} && \Phi r - T\Phi r \leq s + \gamma \mathbf{1}, \\ &&& \pi_{\mu^*, \nu}^\top s \leq \theta, \quad s \geq 0. \end{aligned}$$

Suppose we are given a vector  $r$  of basis function weights and a violation budget  $\theta$ . As we will shortly demonstrate,  $\ell(r, \theta)$  defines the minimum translation (in the direction of the vector  $\mathbf{1}$ ) of  $r$  such so as to get a feasible solution for (8). We will denote by  $s(r, \theta)$  the  $s$  component of the solution to (9). The following lemma, whose proof may be found in Appendix A, characterizes the function  $\ell(r, \theta)$ :

**Lemma 1.** *For any  $r \in \mathbb{R}^K$  and  $\theta \geq 0$ :*

(i)  $\ell(r, \theta)$  is a finite-valued, decreasing, piecewise linear, convex function of  $\theta$ .

(ii)

$$\ell(r, \theta) \leq \frac{1 + \alpha}{1 - \alpha} \|J^* - \Phi r\|_\infty.$$

(iii) The right partial derivative of  $\ell(r, \theta)$  with respect to  $\theta$  satisfies

$$\frac{\partial^+}{\partial \theta^+} \ell(r, 0) = - \left( (1 - \alpha) \sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x) \right)^{-1},$$

where

$$\Omega(r) \triangleq \operatorname{argmax}_{x \in \mathcal{X}} \Phi r(x) - T\Phi r(x).$$

Armed with this definition, we are now in a position to state our first, crude approximation guarantee:

**Theorem 1.** *Suppose that  $r_{\text{SALP}}$  is an optimal solution to the SALP (8), and let  $r^*$  satisfy*

$$r^* \in \operatorname{argmin}_r \|J^* - \Phi r\|_\infty.$$

Then,

$$(10) \quad \|\Phi r_{\text{SALP}} - J^*\|_{1, \nu} \leq \|J^* - \Phi r^*\|_\infty + \ell(r^*, \theta) + \frac{2\theta}{1 - \alpha}.$$

The above theorem allows us to interpret  $\ell(r^*, \theta) + 2\theta/(1 - \alpha)$  as an upper bound to the approximation error (in the  $\|\cdot\|_{1, \nu}$  norm) associated with the SALP solution  $r_{\text{SALP}}$ , relative to the error of the *best* approximation  $r^*$  (in the  $\|\cdot\|_\infty$  norm). This theorem also provides justification for the intuition, described in Section 3, that a relaxation of the feasible region of the ALP will result in better value function approximations. To see this, consider the following corollary:

**Corollary 1.** Define  $U_{SALP}(\theta)$  to be the upper bound in (10), i.e.,

$$U_{SALP}(\theta) \triangleq \|J^* - \Phi r^*\|_\infty + \ell(r^*, \theta) + \frac{2\theta}{1 - \alpha}.$$

Then:

(i)

$$U_{SALP}(0) \leq \frac{2}{1 - \alpha} \|J^* - \Phi r^*\|_\infty.$$

(ii) The right partial derivative of  $U_{SALP}(\theta)$  with respect to  $\theta$  satisfies

$$\frac{d^+}{d\theta^+} U_{SALP}(0) = \frac{1}{1 - \alpha} \left[ 2 - \left( \sum_{x \in \Omega(r^*)} \pi_{\mu^*, \nu}(x) \right)^{-1} \right].$$

**Proof.** The result follows immediately from Parts (ii) and (iii) of Lemma 1. ■

Suppose that  $\theta = 0$ , in which case the SALP (8) is identical to the ALP (3), thus,  $r_{SALP} = r_{ALP}$ . Applying Part (i) of Corollary 1, we have, for the ALP, the approximation error bound

$$(11) \quad \|J^* - \Phi r_{ALP}\|_{1, \nu} \leq \frac{2}{1 - \alpha} \|J^* - \Phi r^*\|_\infty.$$

This is precisely Theorem 2 of de Farias and Van Roy (2003); we recover their approximation guarantee for the ALP.

Now observe that, from Part (ii) of Corollary 1, if the set  $\Omega(r^*)$  is of very small probability according to the distribution  $\pi_{\mu^*, \nu}$ , we expect that the upper bound  $U_{SALP}(\theta)$  will decrease dramatically as  $\theta$  is increased from 0.<sup>2</sup> In other words, if the Bellman error  $\Phi r^*(x) - T\Phi r^*(x)$  produced by  $r^*$  is maximized at states  $x$  that are collectively of small very probability, then we expect to have a choice of  $\theta > 0$  for which

$$U_{SALP}(\theta) \ll U_{SALP}(0) \leq \frac{2}{1 - \alpha} \|J^* - \Phi r^*\|_\infty.$$

In this case, the bound (10) on the SALP solution will be an improvement over the bound (11) on the ALP solution.

Before we present the proof of Theorem 1 we present an auxiliary claim that we will have several opportunities to use. The proof can be found in Appendix A.

**Lemma 2.** Suppose that the vectors  $J \in \mathbb{R}^{\mathcal{X}}$  and  $s \in \mathbb{R}^{\mathcal{X}}$  satisfy

$$J \leq T_{\mu^*} J + s.$$

---

<sup>2</sup>Already if  $\pi_{\mu^*, \nu}(\Omega(r^*)) < 1/2$ ,  $\frac{d^+}{d\theta^+} U_{SALP}(0) < 0$ .

Then,

$$J \leq J^* + \Delta^* s,$$

where

$$\Delta^* \triangleq \sum_{k=0}^{\infty} (\alpha P_{\mu^*})^k = (I - \alpha P_{\mu^*})^{-1},$$

and  $P_{\mu^*}$  is the transition probability matrix corresponding to an optimal policy.

In particular, if  $(r, s)$  is feasible for the LP (8). Then,

$$\Phi r \leq J^* + \Delta^* s.$$

A feasible solution to the ALP is necessarily a lower bound to the optimal cost-to-go function,  $J^*$ . This is no longer the case for the SALP; the above lemma characterizes the extent to which this restriction is relaxed.

We now proceed with the proof of Theorem 1:

**Proof of Theorem 1.** First, define the weight vector  $\tilde{r} \in \mathbb{R}^m$  by

$$\Phi \tilde{r} = \Phi r^* - \ell(r^*, \theta) \mathbf{1},$$

and set  $\tilde{s} = s(r^*, \theta)$ , the  $s$ -component of the solution to the LP (9) with parameters  $r^*$  and  $\theta$ . We will demonstrate that  $(\tilde{r}, \tilde{s})$  is feasible for (5). Observe that, by the definition of the LP (9),

$$\Phi r^* \leq T\Phi r^* + \tilde{s} + (1 - \alpha)\ell(r^*, \theta)\mathbf{1}.$$

Then,

$$\begin{aligned} T\Phi \tilde{r} &= T\Phi r^* - \alpha\ell(r^*, \theta)\mathbf{1} \\ &\geq \Phi r^* - \tilde{s} - (1 - \alpha)\ell(r^*, \theta)\mathbf{1} - \alpha\ell(r^*, \theta)\mathbf{1} \\ &= \Phi \tilde{r} - \tilde{s}. \end{aligned}$$

Now, let  $(r_{\text{SALP}}, \bar{s})$  be the solution to the SALP (8). By Lemma 2,

$$\begin{aligned}
\|J^* - \Phi r_{\text{SALP}}\|_{1,\nu} &\leq \|J^* - \Phi r_{\text{SALP}} + \Delta^* \bar{s}\|_{1,\nu} + \|\Delta^* \bar{s}\|_{1,\nu} \\
&= \nu^\top (J^* - \Phi r_{\text{SALP}} + \Delta^* \bar{s}) + \nu^\top \Delta^* \bar{s} \\
&= \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2\pi_{\mu^*, \nu}^\top \bar{s}}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2\theta}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi \tilde{r}) + \frac{2\theta}{1 - \alpha} \\
&\leq \|J^* - \Phi \tilde{r}\|_\infty + \frac{2\theta}{1 - \alpha} \\
&\leq \|J^* - \Phi r^*\|_\infty + \|\Phi r^* - \Phi \tilde{r}\|_\infty + \frac{2\theta}{1 - \alpha} \\
&= \|J^* - \Phi r^*\|_\infty + \ell(r^*, \theta) + \frac{2\theta}{1 - \alpha},
\end{aligned}$$

as desired. ■

While Theorem 1 reinforces the intuition (shown via Figure 1) that the SALP will permit closer approximations to  $J^*$  than the ALP, the bound leaves room for improvement:

1. The right hand side of our bound measures projection error,  $\|J^* - \Phi r^*\|_\infty$  in the  $\|\cdot\|_\infty$  norm. Since it is unlikely that the basis functions  $\Phi$  will provide a uniformly good approximation over the entire state space, the right hand side of our bound could be quite large.
2. As suggested by (4), the choice of state relevance weights can significantly influence the solution. In particular, it allows us to choose regions of the state space where we would like a better approximation of  $J^*$ . The right hand side of our bound, however, is independent of  $\nu$ .
3. Our guarantee does not suggest a concrete choice of the violation budget parameter  $\theta$ .

The next section will present a substantially refined approximation bound, that will address these issues.

### 4.3. A Stronger Approximation Guarantee

With the intent of deriving stronger approximation guarantees, we begin this section by introducing a ‘nicer’ measure of the quality of approximation afforded by  $\Phi$ . In particular, instead of measuring the approximation error  $J^* - \Phi r^*$  in the  $\|\cdot\|_\infty$  norm as we did for our previous bounds, we will use a weighted max norm defined according to:

$$\|J\|_{\infty, 1/\psi} \triangleq \max_{x \in \mathcal{X}} \frac{|J(x)|}{\psi(x)}.$$

Here,  $\psi: \mathcal{X} \rightarrow [1, \infty)$  is a given ‘weighting’ function. The weighting function  $\psi$  allows us to weight approximation error in a non-uniform fashion across the state space and in this manner potentially ignore approximation quality in regions of the state space that are less relevant. We define  $\Psi$  to be the set of all weighting functions, i.e.,

$$\Psi \triangleq \left\{ \psi \in \mathbb{R}^{\mathcal{X}} : \psi \geq \mathbf{1} \right\}.$$

Given a particular  $\psi \in \Psi$ , we define a scalar

$$\beta(\psi) \triangleq \max_{x,a} \left| \frac{\sum_{x'} P_a(x, x') \psi(x')}{\psi(x)} \right|.$$

One may view  $\beta(\psi)$  as a measure of the ‘stability’ of the system.

In addition to specifying the sampling distribution  $\pi$ , as we did in Section 4.2, we will specify (implicitly) a particular choice of the violation budget  $\theta$ . In particular, we will consider solving the following SALP:

$$(12) \quad \begin{aligned} & \underset{r,s}{\text{maximize}} && \nu^\top \Phi r - \frac{2\pi_{\mu^*, \nu}^\top s}{1-\alpha} \\ & \text{subject to} && \Phi r \leq T\Phi r + s, \quad s \geq 0. \end{aligned}$$

It is clear that (12) is equivalent to (8) for a specific choice of  $\theta$ . We then have:

**Theorem 2.** *If  $r_{SALP}$  is an optimal solution to (12), then*

$$\|J^* - \Phi r_{SALP}\|_{1, \nu} \leq \inf_{r, \psi \in \Psi} \|J^* - \Phi r\|_{\infty, 1/\psi} \left( \nu^\top \psi + \frac{2(\pi_{\mu^*, \nu}^\top \psi + 1)(\alpha\beta(\psi) + 1)}{1 - \alpha} \right).$$

Before presenting a proof for this approximation guarantee, it is worth placing the result in context to understand its implications. For this, we recall a closely related result shown by de Farias and Van Roy (2003) for the ALP. In particular, de Farias and Van Roy (2003) showed that *given* an appropriate weighting function (in their context, a ‘Lyapunov’ function)  $\psi$ , one may solve an ALP, with  $\psi$  in the span of the basis functions  $\Phi$ . The solution  $r_{ALP}$  to such an ALP then satisfies:

$$(13) \quad \|J^* - \Phi r_{ALP}\|_{1, \nu} \leq \inf_r \|J^* - \Phi r\|_{\infty, 1/\psi} \frac{2\nu^\top \psi}{1 - \alpha\beta(\psi)},$$

provided that  $\beta(\psi) < 1/\alpha$ . Selecting an appropriate  $\psi$  in their context is viewed to be an important task for practical performance and often requires a good deal of problem specific analysis; de Farias and Van Roy (2003) identify appropriate  $\psi$  for several queueing models. Note that this is equivalent to identifying a desirable basis function. In contrast, the guarantee we present optimizes over *all possible*  $\psi$  (including those  $\psi$  that do not satisfy the Lyapunov condition  $\beta(\psi) < 1/\alpha$ , and that are not necessarily in the span of  $\Phi$ ).

To make the comparison more precise, let us focus attention on a particular choice of  $\nu$ , namely  $\nu = \pi_{\mu^*} \triangleq \pi_*$ , the stationary distribution induced under an optimal policy  $\mu^*$ . In this case, restricting attention to the set of weighting functions

$$\bar{\Psi} = \{\psi \in \Psi : \alpha\beta(\psi) < 1\},$$

so as to make the two bounds comparable, Theorem 2 guarantees that

$$\begin{aligned} \|J^* - \Phi r_{\text{SALP}}\|_{1,\nu} &\leq \inf_{r, \psi \in \bar{\Psi}} \|J^* - \Phi r\|_{\infty, 1/\psi} \left( \pi_*^\top \psi + \frac{2(\pi_*^\top \psi + 1)(\alpha\beta(\psi) + 1)}{1 - \alpha} \right) \\ &\leq \inf_{r, \psi \in \bar{\Psi}} \|J^* - \Phi r\|_{\infty, 1/\psi} \frac{9\pi_*^\top \psi}{1 - \alpha}. \end{aligned}$$

On the other hand, observing that  $\beta(\psi) \geq 1$  for all  $\psi \in \Psi$ , the right hand side for the ALP bound (13) is at least

$$\inf_r \|J^* - \Phi r\|_{\infty, 1/\psi} \frac{2\pi_*^\top \psi}{1 - \alpha}.$$

Thus, the approximation guarantee of Theorem 2 allows us to view the SALP as *automating* the critical procedure of identifying a good Lyapunov function for a given problem.

**Proof of Theorem 2.** Let  $r \in \mathbb{R}^m$  and  $\psi \in \Psi$  be arbitrary. Define the vectors  $\tilde{\epsilon}, \tilde{s} \in \mathbb{R}^{\mathcal{X}}$  component-wise by

$$\begin{aligned} \tilde{\epsilon}(x) &\triangleq ((\Phi r)(x) - (T\Phi r)(x))^+, \\ \tilde{s}(x) &\triangleq \tilde{\epsilon}(x) \left( 1 - \frac{1}{\psi(x)} \right). \end{aligned}$$

Notice that  $0 \leq \tilde{s} \leq \tilde{\epsilon}$ .

We next make a few observations. First, define  $\tilde{r}$  according to

$$\Phi \tilde{r} = \Phi r - \frac{\|\tilde{\epsilon}\|_{\infty, 1/\psi}}{1 - \alpha} \mathbf{1}.$$

Observe that, by a similar construction to that in Theorem 1,  $(\tilde{r}, \tilde{s})$  is feasible for (12). Also,

$$\|\Phi r - \Phi \tilde{r}\|_{\infty} = \frac{\|\tilde{\epsilon}\|_{\infty, 1/\psi}}{1 - \alpha} \leq \frac{\|T\Phi r - \Phi r\|_{\infty, 1/\psi}}{1 - \alpha}.$$

Furthermore,

$$\begin{aligned}
\pi_{\mu^*, \nu}^\top \tilde{s} &= \sum_{x \in \mathcal{X}} \pi_{\mu^*, \nu}(x) \tilde{\epsilon}(x) (1 - 1/\psi(x)) \\
&\leq \pi_{\mu^*, \nu}^\top \tilde{\epsilon} \\
&\leq (\pi_{\mu^*, \nu}^\top \psi) \|\tilde{\epsilon}\|_{\infty, 1/\psi} \\
&\leq (\pi_{\mu^*, \nu}^\top \psi) \|T\Phi r - \Phi r\|_{\infty, 1/\psi}.
\end{aligned}$$

Finally, note that

$$\nu^\top (J^* - \Phi r) \leq (\nu^\top \psi) \|J^* - \Phi r\|_{\infty, 1/\psi}.$$

Now, suppose that  $(r_{\text{SALP}}, \bar{s})$  is an optimal solution to the SALP (12). We have from the last set of inequalities in the proof of Theorem 1 and the above observations,

$$\begin{aligned}
\|J^* - \Phi r_{\text{SALP}}\|_{1, \nu} &\leq \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2\pi_{\mu^*, \nu}^\top \bar{s}}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi \tilde{r}) + \frac{2\pi_{\mu^*, \nu}^\top \tilde{s}}{1 - \alpha} \\
(14) \quad &\leq \nu^\top (J^* - \Phi r) + \nu^\top (\Phi r - \Phi \tilde{r}) + \frac{2\pi_{\mu^*, \nu}^\top \tilde{s}}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi r) + \|\Phi r - \Phi \tilde{r}\|_\infty + \frac{2\pi_{\mu^*, \nu}^\top \tilde{s}}{1 - \alpha} \\
&\leq (\nu^\top \psi) \|J^* - \Phi r\|_{\infty, 1/\psi} + \frac{\|T\Phi r - \Phi r\|_{\infty, 1/\psi}}{1 - \alpha} (1 + 2\pi_{\mu^*, \nu}^\top \psi).
\end{aligned}$$

Since our choice of  $r$  and  $\psi$  were arbitrary, we have:

$$(15) \quad \|J^* - \Phi r_{\text{SALP}}\|_{1, \nu} \leq \inf_{r, \psi \in \Psi} (\nu^\top \psi) \|J^* - \Phi r\|_{\infty, 1/\psi} + \frac{\|T\Phi r - \Phi r\|_{\infty, 1/\psi}}{1 - \alpha} (1 + 2\pi_{\mu^*, \nu}^\top \psi).$$

We would like to relate the Bellman error term  $T\Phi r - \Phi r$  on the right hand side of (15) to the approximation error  $J^* - \Phi r$ . In order to do so, note that for any vectors  $J, \bar{J} \in \mathbb{R}^{\mathcal{X}}$ ,

$$|TJ(x) - T\bar{J}(x)| \leq \alpha \max_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}} P_a(x, x') |J(x') - \bar{J}(x')|.$$

Therefore,

$$\begin{aligned}
\|T\Phi r - J^*\|_{\infty, 1/\psi} &\leq \alpha \max_{x, a} \frac{\sum_{x'} P_a(x, x') |\Phi r(x') - J^*(x')|}{\psi(x)} \\
&\leq \alpha \max_{x, a} \frac{\sum_{x'} P_a(x, x') \psi(x') \frac{|\Phi r(x') - J^*(x')|}{\psi(x')}}{\psi(x)} \\
&\leq \alpha \beta(\psi) \|J^* - \Phi r\|_{\infty, 1/\psi}.
\end{aligned}$$



Thus,

$$(16) \quad \begin{aligned} \|T\Phi r - \Phi r\|_{\infty, 1/\psi} &\leq \|T\Phi r - J^*\|_{\infty, 1/\psi} + \|J^* - \Phi r\|_{\infty, 1/\psi} \\ &\leq \|J^* - \Phi r\|_{\infty, 1/\psi} (1 + \alpha\beta(\psi)). \end{aligned}$$

Combining (15) and (16), we get the desired result. ■

The analytical results provided in Sections 4.2 and 4.3 provide bounds on the quality of the approximation provided by the SALP solution to  $J^*$ . The next section presents performance bounds with the intent of understanding the increase in expected cost incurred in using a control policy that is greedy with respect to the SALP approximation in lieu of the optimal policy.

#### 4.4. A Performance Bound

We will momentarily present a result that will allow us to interpret the objective of the SALP (12) as an upper bound on the performance loss of a greedy policy with respect to the SALP solution. Before doing so, we briefly introduce some relevant notation. For a given policy  $\mu$ , we denote

$$\Delta_\mu \triangleq \sum_{k=0}^{\infty} (\alpha P_\mu)^k = (I - \alpha P_\mu)^{-1}.$$

Thus,  $\Delta^* = \Delta_{\mu^*}$ . Given a vector  $J \in \mathbb{R}^{\mathcal{X}}$ , let  $\mu_J$  denote the greedy policy with respect to  $J$ . That is,  $\mu_J$  satisfies  $T_{\mu_J} J = TJ$ . Recall that the policy of interest to us will be  $\mu_{\Phi r_{\text{SALP}}}$  for a solution  $r_{\text{SALP}}$  to the SALP. Finally, for an arbitrary starting distribution over states  $\eta$ , we define the ‘discounted’ steady state distribution over states induced by  $\mu_J$  according to

$$\nu(\eta, J)^\top \triangleq (1 - \alpha)\eta^\top \sum_{k=0}^{\infty} (\alpha P_{\mu_J})^k = (1 - \alpha)\eta^\top \Delta_{\mu_J}.$$

We have the following upper bound on the increase in cost incurred by using  $\mu_J$  in place of  $\mu^*$ :

**Theorem 3.**

$$\|J_{\mu_J} - J^*\|_{1, \eta} \leq \frac{1}{1 - \alpha} \left( \nu(\eta, J)^\top (J^* - J) + \frac{2}{1 - \alpha} \pi_{\mu^*, \nu(\eta, J)}^\top (J - TJ)^+ \right).$$

Theorem 3 indicates that if  $J$  is close to  $J^*$ , so that  $(J - TJ)^+$  is also small, then the expected cost incurred in using a control policy that is greedy with respect to  $J$  will be close to optimal. The bound indicates the impact of approximation errors over differing parts of the state space on performance loss.

Suppose that  $(r_{\text{SALP}}, \bar{s})$  is an optimal solution to the SALP (12). Then, examining the proof of

Theorem 2 and, in particular, (14), reveals that

$$\begin{aligned}
(17) \quad & \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2}{1-\alpha} \pi_{\mu^*, \nu}^\top \bar{s} \\
& \leq \inf_{r, \psi \in \Psi} \|J^* - \Phi r\|_{\infty, 1/\psi} \left( \nu^\top \psi + \frac{2(\pi_{\mu^*, \nu}^\top \psi + 1)(\alpha\beta(\psi) + 1)}{1-\alpha} \right).
\end{aligned}$$

Assume that the state relevance weights  $\nu$  in the SALP (12) satisfy

$$(18) \quad \nu = \nu(\eta, \Phi r_{\text{SALP}}).$$

Then, combining Theorem 2 and (17) yields

$$(19) \quad \|J_{\mu\Phi r_{\text{SALP}}} - J^*\|_{1, \eta} \leq \frac{1}{1-\alpha} \left( \inf_{r, \psi \in \Psi} \|J^* - \Phi r\|_{\infty, 1/\psi} \left( \nu^\top \psi + \frac{2(\pi_{\mu^*, \nu}^\top \psi + 1)(\alpha\beta(\psi) + 1)}{1-\alpha} \right) \right).$$

This bound *directly* relates the performance loss of the SALP policy to the ability of the basis function architecture  $\Phi$  to approximate  $J^*$ . Moreover, this relationship allows us to loosely interpret the SALP as minimizing an upper bound on performance loss.

Unfortunately, it is not clear how to make an a-priori choice of the state relevance weights  $\nu$  to satisfy (18), since the choice of  $\nu$  determines the solution to the SALP  $r_{\text{SALP}}$ ; this is essentially the situation one faces in performance analyses for approximate dynamic programming algorithms such as approximate value iteration and temporal difference learning (de Farias and Van Roy, 2000). Indeed, it is not clear that there exists a  $\nu$  that solves the fixed point equation (18). On the other hand, given a choice of  $\nu$  so that  $\nu \approx \nu(\eta, \Phi r_{\text{SALP}})$ , in the sense of a small Radon-Nikodym derivative between the two distributions, an approximate version of the performance bound (19) will hold. As suggested by de Farias and Van Roy (2003) in the ALP case, one possibility for finding such a choice of state relevance weights is to iteratively resolve the SALP, and at each time using the policy from the prior iteration to generate state relevance weights for the next iteration.

**Proof of Theorem 3.** Define  $s \triangleq (J - TJ)^+$ . From Lemma 2, we know that

$$J \leq J^* + \Delta^* s.$$

Applying  $T_{\mu^*}$  to both sides,

$$T_{\mu^*} J \leq J^* + \alpha P_{\mu^*} \Delta^* s = J^* + \Delta^* s - s \leq J^* + \Delta^* s,$$

so that

$$(20) \quad TJ \leq T_{\mu^*} J \leq J^* + \Delta^* s.$$

Then,

$$\begin{aligned}
(21) \quad \eta^\top (J_{\mu_J} - J) &= \eta^\top \sum_{k=0}^{\infty} \alpha^k P_{\mu_J}^k (g_\mu + \alpha P_{\mu_J} J - J) \\
&= \eta^\top \Delta_{\mu_J} (TJ - J) \\
&\leq \eta^\top \Delta_{\mu_J} (J^* - J + \Delta^* s) \\
&= \frac{1}{1-\alpha} \nu(\eta, J)^\top (J^* - J + \Delta^* s).
\end{aligned}$$

where the second equality is from the fact that  $g_\mu + \alpha P_{\mu_J} J = T_{\mu_J} J = TJ$ , and the inequality follows from (20).

Further,

$$\begin{aligned}
(22) \quad \eta^\top (J - J^*) &\leq \eta^\top \Delta^* s \\
&\leq \eta^\top \Delta_{\mu_J} \Delta^* s \\
&= \frac{1}{1-\alpha} \nu(\eta, J)^\top \Delta^* s.
\end{aligned}$$

where the second inequality follows from the fact that  $\Delta^* s \geq 0$  and  $\Delta_{\mu_J} = I + \sum_{k=1}^{\infty} \alpha^k P_{\mu_J}^k$ .

It follows from (21) and (22) that

$$\begin{aligned}
\eta^\top (J_{\mu_J} - J^*) &= \eta^\top (J_{\mu_J} - J) + \eta^\top (J - J^*) \\
&\leq \frac{1}{1-\alpha} \nu(\eta, J)^\top (J^* - J + 2\Delta^* s) \\
&= \frac{1}{1-\alpha} \left( \nu(\eta, J)^\top (J^* - J) + \frac{2}{1-\alpha} \pi_{\mu^*, \nu(\eta, J)}^\top s \right),
\end{aligned}$$

which is the result. ■

#### 4.5. Sample Complexity

Our analysis thus far has assumed we have the ability to solve the SALP. The number of constraints and variables in the SALP grows linearly with the size of the state space  $\mathcal{X}$ . Hence, this program will typically be intractable for problems of interest. One solution, which we describe here, is to consider a *sampled* variation of the SALP, where states and constraints are sampled rather than exhaustively considered. In this section, we will argue that the solution to the SALP is well approximated by the solution to a tractable, sampled variation.

In particular, let  $\hat{\mathcal{X}}$  be a collection of  $S$  states drawn independently from the state space  $\mathcal{X}$

according to the distribution  $\pi_{\mu^*, \nu}$ . Consider the following optimization program:

$$\begin{aligned}
(23) \quad & \underset{r, s}{\text{maximize}} && \nu^\top \Phi r - \frac{2}{(1-\alpha)S} \sum_{x \in \hat{\mathcal{X}}} s(x) \\
& \text{subject to} && \Phi r(x) \leq T\Phi r(x) + s(x), \quad \forall x \in \hat{\mathcal{X}}, \\
& && s \geq 0, \quad r \in \mathcal{N}.
\end{aligned}$$

Here,  $\mathcal{N} \subset \mathbb{R}^K$  is a bounding set that restricts the magnitude of the sampled SALP solution, we will discuss the role of  $\mathcal{N}$  shortly. Notice that (23) is a variation of (12), where only the decision variables and constraints corresponding to the sampled subset of states are retained. The resulting optimization program has  $K + S$  decision variables and  $S|\mathcal{A}|$  linear constraints. For a moderate number of samples  $S$ , this is easily solved. Even in scenarios where the size of the action space  $\mathcal{A}$  is large, it is frequently possible to rewrite (23) as a compact linear program (Farias and Van Roy, 2007; Moallemi et al., 2008). The natural question, however, is whether the solution to the sampled SALP (23) is a good approximation to the solution provided by the SALP (12), for a ‘tractable’ number of samples  $S$ .

Here, we answer this question in the affirmative. We will provide a sample complexity bound that indicates that for a number of samples  $S$  that scales linearly with the dimension of  $\Phi$ ,  $K$ , and that need not depend on the size of the state space, the solution to the sampled SALP satisfies, with high probability, the approximation guarantee presented for the SALP solution in Theorem 2.

Our proof will rely on the following lemma, which provides a Chernoff bound for the *uniform* convergence of a certain class of functions. The proof of this lemma, which is based on bounding the pseudo-dimension of the class of functions, can be found in Appendix B.

**Lemma 3.** *Given a constant  $B > 0$ , define the function  $\zeta: \mathbb{R} \rightarrow [0, B]$  by*

$$\zeta(t) \triangleq \max(\min(t, B), 0).$$

*Consider a pair of random variables  $(Y, Z) \in \mathbb{R}^K \times \mathbb{R}$ . For each  $i = 1, \dots, n$ , let the pair  $(Y^{(i)}, Z^{(i)})$  be an i.i.d. sample drawn according to the distribution of  $(Y, Z)$ . Then, for all  $\epsilon \in (0, B]$ ,*

$$\begin{aligned}
\mathbb{P} \left( \sup_{r \in \mathbb{R}^K} \left| \frac{1}{n} \sum_{i=1}^n \zeta(r^\top Y^{(i)} + Z^{(i)}) - \mathbb{E} [\zeta(r^\top Y + Z)] \right| > \epsilon \right) \\
\leq 8 \left( \frac{32eB}{\epsilon} \log \frac{32eB}{\epsilon} \right)^{K+2} \exp \left( -\frac{\epsilon^2 n}{64B^2} \right).
\end{aligned}$$

Moreover, given  $\delta \in (0, 1)$ , if

$$n \geq \frac{64B^2}{\epsilon^2} \left( 2(K+2) \log \frac{16eB}{\epsilon} + \log \frac{8}{\delta} \right),$$

then this probability is at most  $\delta$ .

In order to establish a sample complexity result, we require control over the magnitude of optimal solutions to the SALP (12). This control is provided by the bounding set  $\mathcal{N}$ . In particular, we will assume that  $\mathcal{N}$  is large enough so that it contains an optimal solution to the SALP (12), and we define the constant

$$(24) \quad B \triangleq \sup_{r \in \mathcal{N}} \|(\Phi r - T\Phi r)^+\|_\infty.$$

This quantity is closely related to the diameter of the region  $\mathcal{N}$ . Our main sample complexity result can then be stated as follows:

**Theorem 4.** *Under the conditions of Theorem 2, let  $r_{SALP}$  be an optimal solution to the SALP (12), and let  $\hat{r}_{SALP}$  be an optimal solution to the sampled SALP (23). Assume that  $r_{SALP} \in \mathcal{N}$ . Further, given  $\epsilon \in (0, B]$  and  $\delta \in (0, 1/2]$ , suppose that the number of sampled states  $S$  satisfies*

$$S \geq \frac{64B^2}{\epsilon^2} \left( 2(K+2) \log \frac{16eB}{\epsilon} + \log \frac{8}{\delta} \right).$$

Then, with probability at least  $1 - \delta - 2^{-383}\delta^{128}$ ,

$$\|J^* - \Phi \hat{r}_{SALP}\|_{1,\nu} \leq \inf_{\substack{r \in \mathcal{N} \\ \psi \in \Psi}} \|J^* - \Phi r\|_{\infty, 1/\psi} \left( \nu^\top \psi + \frac{2(\pi_{\mu^*, \nu}^\top \psi + 1)(\alpha\beta(\psi) + 1)}{1 - \alpha} \right) + \frac{4\epsilon}{1 - \alpha}.$$

Theorem 4 establishes that the sampled SALP (23) provides a close approximation to the solution of the SALP (12), in the sense that the approximation guarantee we established for the SALP in Theorem 2 is approximately valid for the solution to the sampled SALP, with high probability. The theorem precisely specifies the number of samples required to accomplish this task. This number depends linearly on the number of basis functions and the diameter of the feasible region, but is otherwise independent of the size of the state space for the MDP under consideration.

It is worth juxtaposing our sample complexity result with that available for the ALP (3). Recall that the ALP has a large number of constraints but a *small* number of variables; the SALP is thus, at least superficially, a significantly more complex program. Exploiting the fact that the ALP has a small number of variables, de Farias and Van Roy (2004) establish a sample complexity bound for a sampled version of the ALP analogous to the the sampled SALP (23). The number of samples required for this sampled ALP to produce a good approximation to the ALP can be shown to depend on the same problem parameters we have identified here, viz.: the constant  $B$  and the number of basis functions  $K$ . The sample complexity in the ALP case is identical to the sample complexity bound established here, up to constants and a linear dependence on the ratio  $B/\epsilon$ . This is as opposed to the quadratic dependence on  $B/\epsilon$  of the sampled SALP. Although the two sample complexity bounds are within polynomial terms of each other, one may rightfully worry about the

practical implications of an additional factor of  $B/\epsilon$  in the required number of samples. In the computational study of Section 6, we will attempt to address this concern.

Finally, note that the sampled SALP has  $K + S$  variables and  $S|\mathcal{A}|$  linear constraints whereas the sampled ALP has merely  $K$  variables and  $S|\mathcal{A}|$  linear constraints. Nonetheless, we will show in the Section 5.1 that the special structure of the Hessian associated with the sampled SALP affords us a linear computational complexity dependence on  $S$ .

**Proof of Theorem 4.** Define the vectors

$$\hat{s}_{\mu^*} \triangleq (\Phi \hat{r}_{\text{SALP}} - T_{\mu^*} \Phi \hat{r}_{\text{SALP}})^+, \quad \text{and} \quad \hat{s} \triangleq (\Phi \hat{r}_{\text{SALP}} - T \Phi \hat{r}_{\text{SALP}})^+.$$

One has, via Lemma 2, that

$$\Phi \hat{r}_{\text{SALP}} - J^* \leq \Delta^* \hat{s}_{\mu^*}$$

Thus, as in the last set of inequalities in the proof of Theorem 1, we have

$$(25) \quad \|J^* - \Phi \hat{r}_{\text{SALP}}\|_{1,\nu} \leq \nu^\top (J^* - \Phi \hat{r}_{\text{SALP}}) + \frac{2\pi_{\mu^*,\nu}^\top \hat{s}_{\mu^*}}{1 - \alpha}.$$

Now, let  $\hat{\pi}_{\mu^*,\nu}$  be the empirical measure induced by the collection of sampled states  $\hat{\mathcal{X}}$ . Given a state  $x \in \mathcal{X}$ , define a vector  $Y(x) \in \mathbb{R}^K$  and a scalar  $Z(x) \in \mathbb{R}$  according to

$$Y(x) \triangleq \Phi(x)^\top - \alpha P_{\mu^*} \Phi(x)^\top, \quad Z(x) \triangleq -g(x, \mu^*(x)),$$

so that, for any vector of weights  $r \in \mathcal{N}$ ,

$$(\Phi r(x) - T_{\mu^*} \Phi r(x))^+ = \zeta \left( r^\top Y(x) + Z(x) \right).$$

Then,

$$\left| \hat{\pi}_{\mu^*,\nu}^\top \hat{s}_{\mu^*} - \pi_{\mu^*,\nu}^\top \hat{s}_{\mu^*} \right| \leq \sup_{r \in \mathcal{N}} \left| \frac{1}{S} \sum_{x \in \hat{\mathcal{X}}} \zeta \left( r^\top Y(x) + Z(x) \right) - \sum_{x \in \mathcal{X}} \pi_{\mu^*,\nu}(x) \zeta \left( r^\top Y(x) + Z(x) \right) \right|.$$

Applying Lemma 3, we have that

$$(26) \quad \mathbb{P} \left( \left| \hat{\pi}_{\mu^*,\nu}^\top \hat{s}_{\mu^*} - \pi_{\mu^*,\nu}^\top \hat{s}_{\mu^*} \right| > \epsilon \right) \leq \delta.$$

Next, suppose  $(r_{\text{SALP}}, \bar{s})$  is an optimal solution to the SALP (12). Then, with probability at

least  $1 - \delta$ ,

$$\begin{aligned}
(27) \quad \nu^\top (J^* - \Phi \hat{r}_{\text{SALP}}) + \frac{2\pi_{\mu^*, \nu}^\top \hat{s}_{\mu^*}}{1 - \alpha} &\leq \nu^\top (J^* - \Phi \hat{r}_{\text{SALP}}) + \frac{2\hat{\pi}_{\mu^*, \nu}^\top \hat{s}_{\mu^*}}{1 - \alpha} + \frac{2\epsilon}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi \hat{r}_{\text{SALP}}) + \frac{2\hat{\pi}_{\mu^*, \nu}^\top \hat{s}}{1 - \alpha} + \frac{2\epsilon}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2\hat{\pi}_{\mu^*, \nu}^\top \bar{s}}{1 - \alpha} + \frac{2\epsilon}{1 - \alpha},
\end{aligned}$$

where the first inequality follows from (26), and the final inequality follows from the optimality of  $(\hat{r}_{\text{SALP}}, \hat{s})$  for the sampled SALP (23).

Notice that, without loss of generality, we can assume that  $\bar{s}(x) = (\Phi r_{\text{SALP}}(x) - T\Phi r_{\text{SALP}}(x))^+$ , for each  $x \in \mathcal{X}$ . Thus,  $0 \leq \bar{s}(x) \leq B$ . Applying Hoeffding's inequality,

$$\mathbb{P} \left( \left| \hat{\pi}_{\mu^*, \nu}^\top \bar{s} - \pi_{\mu^*, \nu}^\top \bar{s} \right| \geq \epsilon \right) \leq 2 \exp \left( -\frac{2S\epsilon^2}{B^2} \right) < 2^{-383} \delta^{128},$$

where final inequality follows from our choice of  $S$ . Combining this with (25) and (27), with probability at least  $1 - \delta - 2^{-383} \delta^{128}$ , we have

$$\begin{aligned}
\|J^* - \Phi \hat{r}_{\text{SALP}}\|_{1, \nu} &\leq \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2\hat{\pi}_{\mu^*, \nu}^\top \bar{s}}{1 - \alpha} + \frac{2\epsilon}{1 - \alpha} \\
&\leq \nu^\top (J^* - \Phi r_{\text{SALP}}) + \frac{2\pi_{\mu^*, \nu}^\top \bar{s}}{1 - \alpha} + \frac{4\epsilon}{1 - \alpha}.
\end{aligned}$$

The result then follows from (14)–(16) in the proof of Theorem 2. ■

An alternative sample complexity bound of a similar flavor can be developed using results from the stochastic programming literature. The key idea is that the SALP (12) can be reformulated as the following convex stochastic programming problem:

$$(28) \quad \underset{r \in \mathcal{N}}{\text{maximize}} \quad \mathbb{E}_{\nu, \pi_{\mu^*, \nu}} \left[ \Phi r(x_0) - \frac{2}{1 - \alpha} (\Phi r(x) - T\Phi r(x))^+ \right],$$

where  $x_0, x \in \mathcal{X}$  have distributions  $\nu$  and  $\pi_{\mu^*, \nu}$ , respectively. Interpreting the sampled SALP (23) as a sample average approximation of (28), a sample complexity bound can be developed using the methodology of Shapiro et al. (2009, Chap. 5), for example. This proof is simpler than the one presented here, but yields a cruder estimate that is not as easily compared with those available for the ALP.

## 5. Practical Implementation

The SALP (5), as it is written, is not directly implementable. As discussed in Section 4.5, the number of variables and constraints grows linearly with the size of the state space  $\mathcal{X}$ , making the optimization problem intractable. Moreover, it is not clear how to choose parameters such as the probability distributions  $\nu$  and  $\pi$  or the violation budget  $\theta$ . However, the analysis in Section 4 provides insight that allows us to codify a recipe for a practical and implementable variation.

Consider the following algorithm:

1. Sample  $S$  states independently from the state space  $\mathcal{X}$  according to a sampling distribution  $\rho$ . Denote the set of sampled states by  $\hat{\mathcal{X}}$ .
2. Perform a line search over increasing choices of  $\theta \geq 0$ . For each choice of  $\theta$ ,
  - (a) Solve the *sampled* SALP:

$$\begin{aligned}
 (29) \quad & \underset{r,s}{\text{maximize}} && \frac{1}{S} \sum_{x \in \hat{\mathcal{X}}} (\Phi r)(x) \\
 & \text{subject to} && \Phi r(x) \leq T\Phi r(x) + s(x), \quad \forall x \in \hat{\mathcal{X}}, \\
 & && \frac{1}{S} \sum_{x \in \hat{\mathcal{X}}} s(x) \leq \theta, \\
 & && s \geq 0.
 \end{aligned}$$

- (b) Evaluate the performance of the policy resulting from (29) via Monte Carlo simulation.
3. Select the best of the policies evaluated in Step 2.

This algorithm takes as inputs the following parameters:

- $\Phi$ , a collection of  $K$  basis functions.
- $S$ , the number of states to sample. By sampling  $S$  states, we limit the number of variables and constraints in the sampled SALP (29). Thus, by keeping  $S$  small, the sampled SALP becomes tractable to solve numerically. On the other hand, the quality of the approximation provided by the sampled SALP may suffer if  $S$  is chosen to be too small. The sample complexity theory developed in Section 4.5 suggests that  $S$  can be chosen to grow linearly with  $K$ , the size of the basis set. In particular, a reasonable choice of  $S$  need not depend on the size of the underlying state space.

In practice, we choose  $S \gg K$  to be as large as possible subject to limits on the CPU time and memory required to solve (29). In Section 5.1, we will discuss how the program (29) can be solved efficiently via barrier methods for large choices of  $S$ .



- $\rho$ , a sampling distribution on the state space  $\mathcal{X}$ . The distribution  $\rho$  is used, via Monte Carlo sampling, in place of both the distributions  $\nu$  and  $\pi$  in the SALP (5). Recall that the bounds in Theorems 1 and 2 provide approximation guarantees in a  $\nu$ -weighted 1-norm. This suggests that  $\nu$  should be chosen to emphasize regions of the state space where the quality of approximation is most important. Similarly, the theory in Section 4 suggests that the distribution  $\pi$  should be related to the distribution induced by the optimal policy.

In practice, we choose  $\rho$  to be the stationary distribution under a baseline policy. States are then sampled from  $\rho$  via Monte Carlo simulation of the baseline policy. This baseline policy can correspond, for example, to a heuristic control policy for the system. More sophisticated procedures such as ‘bootstrapping’ can also be considered (Farias and Van Roy, 2006). Here, one starts with a heuristic policy to be used for sampling states. Given the sampled states, the application of our algorithm will result in a new control policy. The new control policy can then be used for state sampling in a subsequent round of optimization, and the process can be repeated.

Note that our algorithm does not require an explicit choice of the violation budget  $\theta$ , since we optimize with a line search over the choices of  $\theta$ . This is motivated by the fact that the sampled SALP (29) can efficiently be resolved for increasing values of  $\theta$  via a ‘warm-start’ procedure. Here, the optimal solution of the sampled SALP given previous value of  $\theta$  is used as a starting point for the solver in a subsequent round of optimization. Using this method we observe that, in practice, the total solution time for a series of sampled SALP instances that vary by their values of  $\theta$  grows sub-linearly with the number of instances.

### 5.1. Efficient Linear Programming Solution

The sampled SALP (29) can be written explicitly in the form of a linear program:

$$(30) \quad \begin{aligned} & \underset{r,s}{\text{maximize}} && c^\top r \\ & \text{subject to} && \begin{bmatrix} A_{11} & A_{12} \\ 0 & d^\top \end{bmatrix} \begin{bmatrix} r \\ s \end{bmatrix} \leq b, \\ & && s \geq 0. \end{aligned}$$

Here,  $b \in \mathbb{R}^{S|\mathcal{A}|+1}$ ,  $c \in \mathbb{R}^K$ , and  $d \in \mathbb{R}^S$  are vectors,  $A_{11} \in \mathbb{R}^{S|\mathcal{A}| \times K}$  is a dense matrix, and  $A_{12} \in \mathbb{R}^{S|\mathcal{A}| \times S}$  is a sparse matrix. This LP has  $K + S$  decision variables and  $S|\mathcal{A}| + 1$  linear constraints.

Typically, the number of sampled states  $S$  will be quite large. For example, in Section 6, we will discuss an example where  $K = 22$  and  $S = 300,000$ . The resulting LP has approximately 300,000 variables and 6,600,000 constraints. In such cases, with many variables *and* many constraints, one might expect the LP to be difficult to solve. However, the sparsity structure of the constraint matrix in (30) and, especially, that of the sub-matrix  $A_{12}$ , allows efficient optimization of this LP.

In particular, imagine solving the LP (30) with a barrier method. The computational bottleneck of such a method is the inner Newton step to compute a central point (see, for example, Boyd and Vandenberghe, 2004). This step involves the solution of a system of linear equations of the form

$$(31) \quad H \begin{bmatrix} \Delta r \\ \Delta s \end{bmatrix} = -g.$$

Here,  $g \in \mathbb{R}^{K+S}$  is a vector and  $H \in \mathbb{R}^{(K+S) \times (K+S)}$  is the Hessian matrix of the barrier function. Without exploiting the structure of the matrix  $H$ , this linear system can be solved with  $O((K+S)^3)$  floating point operations. For large values of  $S$ , this may be prohibitive.

Fortunately, the Hessian matrix  $H$  can be decomposed according to the block structure

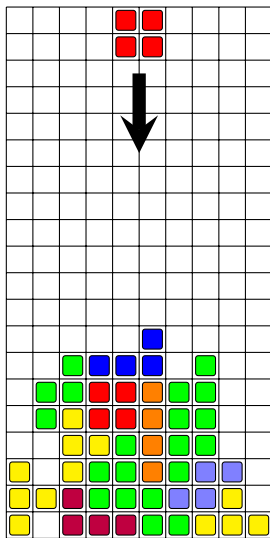
$$H \triangleq \begin{bmatrix} H_{11} & H_{12} \\ H_{12}^\top & H_{22} \end{bmatrix},$$

where  $H_{11} \in \mathbb{R}^{K \times K}$ ,  $H_{12} \in \mathbb{R}^{K \times S}$ , and  $H_{22} \in \mathbb{R}^{S \times S}$ . In the case of the LP (30), it is not difficult to see that the sparsity structure of the sub-matrix  $A_{12}$  ensures that the sub-matrix  $H_{22}$  takes the form of a diagonal matrix plus a rank-one matrix. This allows the linear system (31) to be solved with  $O(K^2S + K^3)$  floating point operations. This is linear in  $S$ , the number of sampled states.

## 6. Case Study: Tetris

Tetris is a popular video game designed and developed by Alexey Pazhitnov in 1985. The Tetris board, illustrated in Figure 2, consists of a two-dimensional grid of 20 rows and 10 columns. The game starts with an empty grid and pieces fall randomly one after another. Each piece consists of four blocks and the player can rotate and translate it in the plane before it touches the ‘floor’. The pieces come in seven different shapes and the next piece to fall is chosen from among these with equal probability. Whenever the pieces are placed such that there is a line of contiguous blocks formed, a point is earned and the line gets cleared. Once the board has enough blocks such that the incoming piece cannot be placed for all translation and rotation, the game terminates. Hence the goal of the player is to clear maximum number of lines before the board gets full.

Our interest in Tetris as a case study for the SALP algorithm is motivated by several facts. First, theoretical results suggest that design of an optimal Tetris player is a difficult problem. Brzustowski (1992) and Burgiel (1997) have shown that the game of Tetris has to end with probability one, under all policies. They demonstrate a sequence of pieces, which leads to termination state of game for all possible actions. Demaine et al. (2003) consider the offline version of Tetris and provide computational complexity results for ‘optimally’ playing Tetris. They show that when the sequence of pieces is known beforehand it is NP-complete to maximize the number of cleared rows, minimize the maximum height of an occupied square, or maximize the number of pieces placed



**Figure 2:** Example of a Tetris board configuration

before the game ends. This suggests that the online version should be computationally difficult.

Second, Tetris represents precisely the kind of large and unstructured MDP for which it is difficult to design heuristic controllers, and hence policies designed by ADP algorithms are particularly relevant. Moreover, Tetris has been employed by a number of researchers as a testbed problem. One of the important steps in applying these techniques is the choice of basis functions. Fortunately, there is a *fixed set of basis functions*, to be described shortly, which have been used by researchers while applying temporal-difference learning (Bertsekas and Ioffe, 1996; Bertsekas and Tsitsiklis, 1996), policy gradient methods (Kakade, 2002), and approximate linear programming (Farias and Van Roy, 2006). Hence, application of SALP to Tetris allows us to make a clear comparison to other ADP methods.

The SALP methodology described in Section 5 was applied as follows:

- **MDP formulation.** We used the formulation of Tetris as a Markov decision problem of Farias and Van Roy (2006). Here, the ‘state’ at a particular time encodes the current board configuration and the shape of the next falling piece, while the ‘action’ determines the placement of the falling piece.
- **Basis functions.** We employed the 22 basis functions originally introduced by Bertsekas and Ioffe (1996). Each basis function takes a Tetris board configuration as its argument. The functions are as follows:
  - Ten basis functions,  $\phi_0, \dots, \phi_9$ , mapping the state to the height  $h_k$  of each of the ten columns.
  - Nine basis functions,  $\phi_{10}, \dots, \phi_{18}$ , each mapping the state to the absolute difference between heights of successive columns:  $|h_{k+1} - h_k|, k = 1, \dots, 9$ .

- One basis function,  $\phi_{19}$ , that maps state to the maximum column height:  $\max_k h_k$
  - One basis function,  $\phi_{20}$ , that maps state to the number of ‘holes’ in the board.
  - One basis function,  $\phi_{21}$ , that is equal to 1 in every state.
- **State sampling.** Given a sample size  $S$ , a collection  $\hat{\mathcal{X}} \subset \mathcal{X}$  of  $S$  states was sampled. These sampled were generated in an i.i.d. fashion from the stationary distribution of a (rather poor) baseline policy<sup>3</sup>. For each choice of sample size  $S$ , ten different collections of  $S$  samples were generated.
  - **Optimization.** Given the collection  $\hat{\mathcal{X}}$  of sampled states, an increasing sequence of choices of the violation budget  $\theta \geq 0$  is considered. For each choice of  $\theta$ , the optimization program (29) was solved.
  - **Policy evaluation.** Given a vector of weights  $\hat{r}$ , the performance of the corresponding policy was evaluated using Monte Carlo simulation. We calculate the average performance of policy  $\mu_{\hat{r}}$  over a series of 3000 games. Performance is measured in terms of the average number of lines eliminated in a single game. The sequence of pieces in each of the 3000 games was fixed across the evaluation of different policies in order to allow better comparisons.

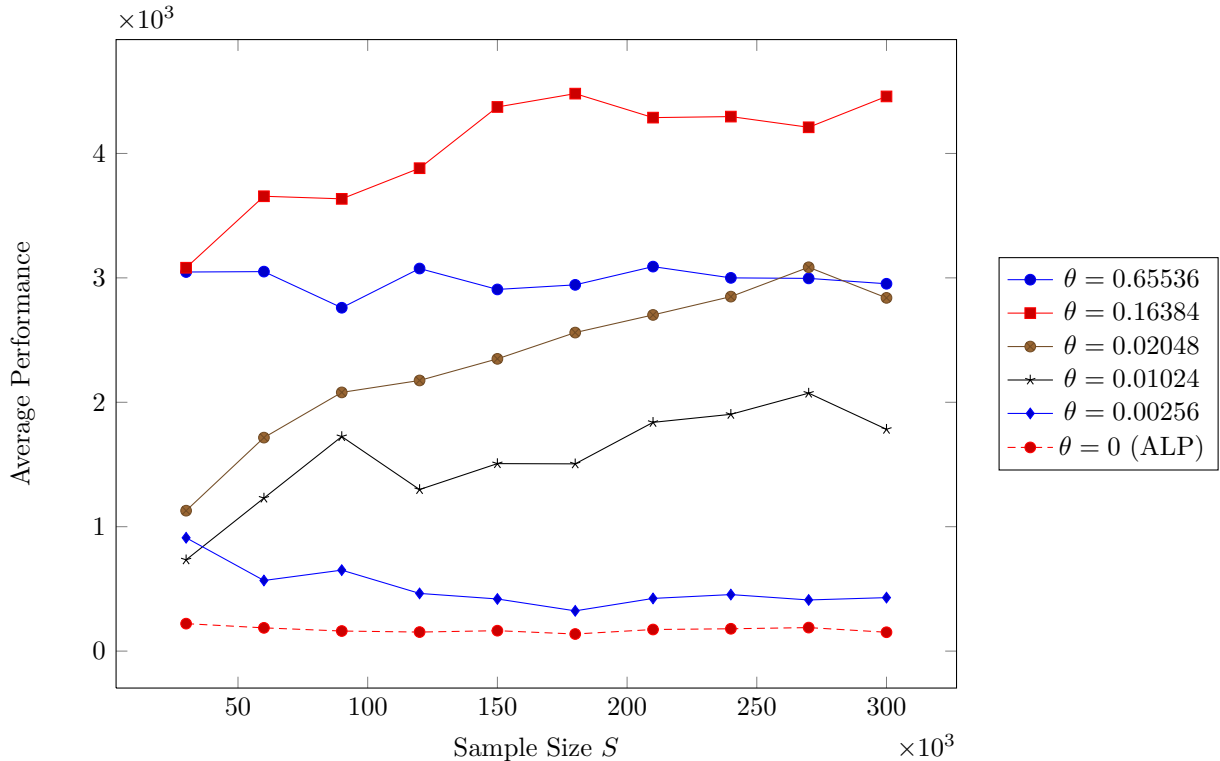
For each pair  $(S, \theta)$ , the resulting *average* performance (averaged over each of the 10 policies arising from the different sets of sampled states) is shown in Figure 3. Note that the  $\theta = 0$  curve in Figure 3 corresponds to the original ALP algorithm. Figure 3 provided experimental evidence for the intuition expressed in Section 3 and the analytic result of Theorem 1: Relaxing the constraints of the ALP even slightly, by allowing for a small slack budget, allows for better policy performance. As the slack budget  $\theta$  is increased from 0, performance dramatically improves. At the peak value of  $\theta = 0.16384$ , the SALP generates policies with performance that is an order of magnitude better than ALP. Beyond this value, the performance of the SALP begins to degrade, as shown by the  $\theta = 0.65536$  curve. Hence, we did not explore larger values of  $\theta$ .

Table 1 summarizes the performance of *best* policies obtained by various ADP algorithms. Note that all of these algorithms employ the same basis function architecture. The ALP and SALP results are from our experiments, while the other results are from the literature. The best performance result of SALP is a factor of 2 better than the competitors.

Note that significantly better policies are possible with this basis function architecture than *any* of the ADP algorithms in Table 1 discover. Using a heuristic global optimization method, Szita and Lőrincz (2006) report finding policies with a remarkable average performance of 350,000. Their method is very computationally intensive, however, requiring one month of CPU time. In addition, the approach employs a number of rather arbitrary Tetris specific ‘modifications’ that are ultimately seen to be critical to performance — in the absence of these modifications, the method is unable to

---

<sup>3</sup>Our baseline policy had an average performance of 113 points.



**Figure 3:** Performance of the average SALP policy for different values of the number of sampled states  $S$  and the violation budget  $\theta$ . Values for  $\theta$  were chosen in an increasing fashion starting from 0, until the resulting average performance began to degrade.

Algorithm	Best Performance	CPU Time
ALP	897	hours
TD-Learning (Bertsekas and Ioffe, 1996)	3,183	minutes
ALP with bootstrapping (Farias and Van Roy, 2006)	4,274	hours
TD-Learning (Bertsekas and Tsitsiklis, 1996)	4,471	minutes
Policy gradient (Kakade, 2002)	5,500	days
SALP	10,775	hours

**Table 1:** Comparison of the performance of the best policy found with various ADP methods.

find a policy for Tetris that scores above a few hundred points. More generally, global optimization methods typically require significant trial and error and other problem specific experimentation in order to work well.

## 7. Conclusion

The approximate linear programming (ALP) approach to approximate DP is interesting at the outset for two reasons. First, the ability to leverage commercial linear programming software to

solve large ADP problems, and second, the ability to prove rigorous approximation guarantees and performance bounds. This paper asked whether the formulation considered in the ALP approach was the ideal formulation. In particular, we asked whether certain strong restrictions imposed on approximations produced by the approach can be relaxed in a tractable fashion and whether such a relaxation has a beneficial impact on the quality of the approximation produced. We have answered both of these questions in the affirmative. In particular, we have presented a novel linear programming formulation that, while remaining no less tractable than the ALP, appears to yield substantial performance gains and permits us to prove extremely strong approximation and performance guarantees.

There are a number of interesting algorithmic directions that warrant exploration. For instance, notice that from (28), that the SALP may be written as an unconstrained stochastic optimization problem. Such problems suggest natural *online* update rules for the weights  $r$ , based on stochastic gradient methods, yielding ‘data-driven’ ADP methods. The menagerie of online ADP algorithms available at present are effectively iterative methods for solving a projected version of Bellman’s equation. TD-learning is a good representative of this type of approach and, as can be seen from Table 1, is not among the highest performing algorithms in our computational study. An online update rule that effectively solves the SALP promises policies that will perform on par with the SALP solution, while at the same time retaining the benefits of an online ADP algorithm. A second interesting algorithmic direction worth exploring is an extension of the smoothed linear programming approach to average cost dynamic programming problems.

As discussed in Section 4, theoretical guarantees for ADP algorithms typically rely on some sort of idealized assumption. For instance, in the case of the ALP, it is the ability to solve an LP with a potentially intractable number of states or else access to a set of sampled states, sampled according to some idealized sampling distribution. For the SALP, it is the latter of the two assumptions. It would be interesting to see whether this assumption can be loosened for some specific class of MDPs. An interesting class of MDPs in this vein are high dimensional optimal stopping problems. Yet another direction for research, is understanding the dynamics of ‘bootstrapping’ procedures, that solve a sequence of sampled versions of the SALP with samples for a given SALP in the sequence drawn according to a policy produced by the previous SALP in the sequence.

## References

- D. Adelman. A price-directed approach to stochastic inventory/routing. *Operations Research*, 52(4):499–514, 2004.
- D. Adelman. Dynamic bid prices in revenue management. *Operations Research*, 55(4):647–661, 2007.
- D. Adelman and D. Klabjan. Computing near optimal policies in generalized joint replenishment. Working paper, January 2009.
- D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 2. Athena Scientific, Belmont, MA, 3rd edition, 2007.

- D. P. Bertsekas and S. Ioffe. Temporal differences-based policy iteration and applications in neuro-dynamic programming. Technical Report LIDS-P-2349, MIT Laboratory for Information and Decision Systems, 1996.
- D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, MA, 1996.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- J. Brzustowski. Can you win at Tetris? Master’s thesis, University of British Columbia, 1992.
- H. Burgiel. How to lose at Tetris. *Mathematical Gazette*, page 194, 1997.
- D. P. de Farias and B. Van Roy. On the existence of fixed points for approximate value iteration and temporal-difference learning. *Journal of Optimization Theory and Applications*, 105(3), 2000.
- D. P. de Farias and B. Van Roy. The linear programming approach to approximate dynamic programming. *Operations Research*, 51(6):850–865, 2003.
- D. P. de Farias and B. Van Roy. On constraint sampling in the linear programming approach to approximate dynamic programming. *Mathematics of Operations Research*, 29(3):462–478, 2004.
- D. P. de Farias and B. Van Roy. A cost-shaping linear program for average-cost approximate dynamic programming with performance guarantees. *Mathematics of Operations Research*, 31(3):597–620, 2006.
- E. D. Demaine, S. Hohenberger, and D. Liben-Nowell. Tetris is hard, even to approximate. In *Proceedings of the 9th International Computing and Combinatorics Conference*, 2003.
- V. F. Farias and B. Van Roy. Tetris: A study of randomized constraint sampling. In *Probabilistic and Randomized Methods for Design Under Uncertainty*. Springer-Verlag, 2006.
- V. F. Farias and B. Van Roy. An approximate dynamic programming approach to network revenue management. Working paper, 2007.
- V. F. Farias, D. Saure, and G. Y. Weintraub. The linear programming approach to solving large scale dynamic stochastic games. Working paper, 2008.
- J. Han. *Dynamic Portfolio Management - An Approximate Linear Programming Approach*. PhD thesis, Stanford University, 2005.
- D. Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 100:78–150, 1992.
- S. Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems 14*, Cambridge, MA, 2002. MIT Press.
- A. S. Manne. Linear programming and sequential decisions. *Management Science*, 60(3):259–267, 1960.
- C. C. Moallemi, S. Kumar, and B. Van Roy. Approximate and data-driven dynamic programming for queueing networks. Working paper, 2008.
- J. R. Morrison and P. R. Kumar. New linear program performance bounds for queueing networks. *Journal of Optimization Theory and Applications*, 100(3):575–597, 1999.
- W. B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. John Wiley and Sons, 2007.
- P. Schweitzer and A. Seidmann. Generalized polynomial approximations in Markovian decision processes. *Journal of Mathematical Analysis and Applications*, 110:568–582, 1985.

- A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, Philadelphia, PA, 2009.
- I. Szita and A. Lőrincz. Learning Tetris using the noisy cross-entropy method. *Neural Computation*, 18: 2936–2941, 2006.
- H. Topaloglu. Using Lagrangian relaxation to compute capacity-dependent bid prices in network revenue management. *Operations Research*, 2009. To appear.
- B. Van Roy. Neuro-dynamic programming: Overview and recent trends. In A. Shwartz E. Feinberg, editor, *Handbook of Markov Decision Processes*. Kluwer, Boston, 2002.
- M. H. Veatch. Approximate dynamic programming for networks: Fluid models and constraint reduction. Working paper, 2005.
- D. Zhang and D. Adelman. An approximate dynamic programming approach to network revenue management with customer choice. Working paper, 2008.

## A. Proofs for Section 4.2

**Lemma 1.** For any  $r \in \mathbb{R}^K$  and  $\theta \geq 0$ :

(i)  $\ell(r, \theta)$  is a finite-valued, decreasing, piecewise linear, convex function of  $\theta$ .

(ii)

$$\ell(r, \theta) \leq \frac{1 + \alpha}{1 - \alpha} \|J^* - \Phi r\|_\infty.$$

(iii) The right partial derivative of  $\ell(r, \theta)$  with respect to  $\theta$  satisfies

$$\frac{\partial^+}{\partial \theta^+} \ell(r, 0) = - \left( (1 - \alpha) \sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x) \right)^{-1},$$

where

$$\Omega(r) \triangleq \operatorname{argmax}_{x \in \mathcal{X}} \Phi r(x) - T\Phi r(x).$$

**Proof.** (i) Given any  $r$ , clearly  $\gamma \triangleq \|\Phi r - T\Phi r\|_\infty$ ,  $s \triangleq 0$  is a feasible point for (9), so  $\ell(r, \theta)$  is well-defined. To see that the LP is bounded, suppose  $(s, \gamma)$  is feasible. Then, for any  $x \in \mathcal{X}$  with  $\pi_{\mu^*, \nu}(x) > 0$ ,

$$\gamma \geq \Phi r(x) - T\Phi r(x) - s(x) \geq \Phi r(x) - T\Phi r(x) - \theta/\pi_{\mu^*, \nu}(x) > \infty.$$

Letting  $(\gamma_1, s_1)$  and  $(\gamma_2, s_2)$  represent optimal solutions for the LP (9) with parameters  $(r, \theta_1)$  and  $(r, \theta_2)$  respectively, it is easy to see that  $((\gamma_1 + \gamma_2)/2, (s_1 + s_2)/2)$  is feasible for the LP with parameters  $(r, (\theta_1 + \theta_2)/2)$ . It follows that  $\ell(r, (\theta_1 + \theta_2)/2) \leq (\ell(r, \theta_1) + \ell(r, \theta_2))/2$ . The remaining properties are simple to check.



(ii) Let  $\epsilon \triangleq \|J^* - \Phi r\|_\infty$ . Then,

$$\|T\Phi r - \Phi r\|_\infty \leq \|J^* - T\Phi r\|_\infty + \|J^* - \Phi r\|_\infty \leq \alpha \|J^* - \Phi r\|_\infty + \epsilon = (1 + \alpha)\epsilon.$$

Since  $\gamma \triangleq \|T\Phi r - \Phi r\|_\infty$ ,  $s \triangleq 0$  is feasible for (9), the result follows.

(iii) Fix  $r \in \mathbb{R}^K$ , and define

$$\Delta \triangleq \max_{x \in \mathcal{X}} (\Phi r(x) - T\Phi r(x)) - \max_{x \in \mathcal{X} \setminus \Omega(r)} (\Phi r(x) - T\Phi r(x)) > 0.$$

Consider the program for  $\ell(r, \delta)$ . It is easy to verify that for  $\delta \geq 0$  and sufficiently small, viz.  $\delta \leq \Delta \sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x)$ ,  $(\bar{s}_\delta, \bar{\gamma}_\delta)$  is an optimal solution to the program, where

$$\bar{s}_\delta(x) \triangleq \begin{cases} \frac{\delta}{\sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x)} & \text{if } x \in \Omega(r), \\ 0 & \text{otherwise,} \end{cases}$$

and

$$\bar{\gamma}_\delta \triangleq \gamma_0 - \frac{\delta}{\sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x)},$$

so that

$$\ell(r, \delta) = \ell(r, 0) - \frac{\delta}{(1 - \alpha) \sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x)}.$$

Thus,

$$\frac{\ell(r, \delta) - \ell(r, 0)}{\delta} = - \left( (1 - \alpha) \sum_{x \in \Omega(r)} \pi_{\mu^*, \nu}(x) \right)^{-1}.$$

Taking a limit as  $\delta \searrow 0$  yields the result. ■

**Lemma 2.** Suppose that the vectors  $J \in \mathbb{R}^{\mathcal{X}}$  and  $s \in \mathbb{R}^{\mathcal{X}}$  satisfy

$$J \leq T_{\mu^*} J + s.$$

Then,

$$J \leq J^* + \Delta^* s,$$

where

$$\Delta^* \triangleq \sum_{k=0}^{\infty} (\alpha P_{\mu^*})^k = (I - \alpha P_{\mu^*})^{-1},$$

and  $P_{\mu^*}$  is the transition probability matrix corresponding to an optimal policy.

In particular, if  $(r, s)$  is feasible for the LP (8). Then,

$$\Phi r \leq J^* + \Delta^* s.$$

**Proof.** Note that the  $T_{\mu^*}$ , the Bellman operator corresponding to the optimal policy  $\mu^*$ , is monotonic and is a contraction. Then, repeatedly applying  $T_{\mu^*}$  to the inequality  $J \leq T_{\mu^*}J + s$  and using the fact that  $T_{\mu^*}^k J \rightarrow J^*$ , we obtain

$$J \leq J^* + \sum_{k=0}^{\infty} (\alpha P_{\mu^*})^k s = J^* + \Delta^* s.$$

■

## B. Proof of Lemma 3

We begin with the following definition: consider a family  $\mathcal{F}$  of functions from a set  $\mathcal{S}$  to  $\{0, 1\}$ . Define the *Vapnik-Chervonenkis (VC) dimension*  $\dim_{VC}(\mathcal{F})$  to be the cardinality  $d$  of the largest set  $\{x_1, x_2, \dots, x_d\} \subset \mathcal{S}$  satisfying:

$$\forall e \in \{0, 1\}^d, \exists f \in \mathcal{F} \text{ such that } \forall i, f(x_i) = 1 \text{ iff } e_i = 1.$$

Now, let  $\mathcal{F}$  be some set of *real*-valued functions mapping  $\mathcal{S}$  to  $[0, B]$ . The *pseudo-dimension*  $\dim_P(\mathcal{F})$  is the following generalization of VC dimension: for each function  $f \in \mathcal{F}$  and scalar  $c \in \mathbb{R}$ , define a function  $g: \mathcal{S} \times \mathbb{R} \rightarrow \{0, 1\}$  according to:

$$g(x, c) \triangleq \mathbb{I}_{\{f(x) - c \geq 0\}}.$$

Let  $\mathcal{G}$  denote the set of all such functions. Then, we define  $\dim_P(\mathcal{F}) \triangleq \dim_{VC}(\mathcal{G})$ .

In order to prove Lemma 3, define the  $\mathcal{F}$  to be the set of functions  $f: \mathbb{R}^K \times \mathbb{R} \rightarrow [0, B]$ , where, for all  $x \in \mathbb{R}^K$  and  $y \in \mathbb{R}$ ,

$$f(y, z) \triangleq \zeta(r^\top y + z).$$

Here,  $\zeta(t) \triangleq \max(\min(t, B), 0)$ , and  $r \in \mathbb{R}^K$  is a vector that parameterizes  $f$ . We will show that  $\dim_P(\mathcal{F}) \leq K + 2$ .

We will use the following standard result from convex geometry:

**Lemma 4 (Radon's Lemma).** *A set  $A \subset \mathbb{R}^m$  of  $m + 2$  points can be partitioned into two disjoint sets  $A_1$  and  $A_2$ , such that the convex hulls of  $A_1$  and  $A_2$  intersect.*

**Lemma 5.**  $\dim_P(\mathcal{F}) \leq K + 2$

**Proof.** Assume, for the sake of contradiction, that  $\dim_P(\mathcal{F}) > K + 2$ . It must be that there exists a ‘shattered’ set

$$\left\{ (y^{(1)}, z^{(1)}, c^{(1)}), (y^{(2)}, z^{(2)}, c^{(2)}), \dots, (y^{(K+3)}, z^{(K+3)}, c^{(K+3)}) \right\} \subset \mathbb{R}^K \times \mathbb{R} \times \mathbb{R},$$

such that, for all  $e \in \{0, 1\}^{K+3}$ , there exists a vector  $r_e \in \mathbb{R}^K$  with

$$\zeta \left( r_e^\top y^{(i)} + z^{(i)} \right) \geq c^{(i)} \text{ iff } e_i = 1, \quad \forall 1 \leq i \leq K+3.$$

Observe that we must have  $c^{(i)} \in (0, B]$  for all  $i$ , since if  $c^{(i)} \leq 0$  or  $c^{(i)} > B$ , then no such shattered set can be demonstrated. But if  $c^{(i)} \in (0, B]$ , for all  $r \in \mathbb{R}^K$ ,

$$\zeta \left( r^\top y^{(i)} + z^{(i)} \right) \geq c^{(i)} \implies r^\top y^{(i)} \geq c^{(i)} - z^{(i)},$$

and

$$\zeta \left( r^\top y^{(i)} + z^{(i)} \right) < c^{(i)} \implies r^\top y^{(i)} < c^{(i)} - z^{(i)}.$$

For each  $1 \leq i \leq K+3$ , define  $x^{(i)} \in \mathbb{R}^{K+1}$  component-wise according to

$$x_j^{(i)} \triangleq \begin{cases} y_j^{(i)} & \text{if } j < K+1, \\ c^{(i)} - z^{(i)} & \text{if } j = K+1. \end{cases}$$

Let  $A = \{x^{(1)}, x^{(2)}, \dots, x^{(K+3)}\} \subset \mathbb{R}^{K+1}$ , and let  $A_1$  and  $A_2$  be subsets of  $A$  satisfying the conditions of Radon's lemma. Define a vector  $\tilde{e} \in \{0, 1\}^{K+3}$  component-wise according to

$$\tilde{e}_i \triangleq \mathbb{I}_{\{x^{(i)} \in A_1\}}.$$

Define the vector  $\tilde{r} \triangleq r_{\tilde{e}}$ . Then, we have

$$\sum_{j=1}^K \tilde{r}_j x_j \geq x_{K+1}, \quad \forall x \in A_1,$$

$$\sum_{j=1}^K \tilde{r}_j x_j < x_{K+1}, \quad \forall x \in A_2.$$

Now, let  $\bar{x} \in \mathbb{R}^{K+1}$  be a point contained in both the convex hull of  $A_1$  and the convex hull of  $A_2$ . Such a point must exist by Radon's lemma. By virtue of being contained in the convex hull of  $A_1$ , we must have

$$\sum_{j=1}^K \tilde{r}_j \bar{x}_j \geq \bar{x}_{K+1}.$$

Yet, by virtue of being contained in the convex hull of  $A_2$ , we must have

$$\sum_{j=1}^K \tilde{r}_j \bar{x}_j < \bar{x}_{K+1},$$

which is impossible. ■

With the above pseudo-dimension estimate, Lemma 3 follows immediately from Corollary 2 of Haussler (1992, Section 4).